

Speaker age estimation using i-vectors

Mohamad Hasan Bahari^{a,*}, Mitchell McLaren^b, Hugo Van hamme^a, David A. van Leeuwen^b

^a*Center for processing speech and images, KU Leuven, Belgium*

^b*Center for Language and Speech Technology, Radboud University Nijmegen, The Netherlands*

Abstract

In this paper, a new approach for age estimation from speech signals based on i-vectors is proposed. In this method, each utterance is modeled by its corresponding i-vector. Then, a Within-Class Covariance Normalization technique is used for session variability compensation. Finally, a least squares support vector regression (LSSVR) is applied to estimate the age of speakers. The proposed method is trained and tested on telephone conversations of the National Institute for Standard and Technology (NIST) 2010 and 2008 speaker recognition evaluation databases. Evaluation results show that the proposed method yields significantly lower mean absolute error and higher Pearson correlation coefficient between chronological speaker age and estimated speaker age compared to different conventional schemes. The obtained relative improvements of mean absolute error and correlation coefficient compared to our best baseline system are around 5% and 2% respectively. Finally, the

*Corresponding author. Tel:+32-(0)16-32.85.45. Fax:+32-(0)16-32.17.23.

Email addresses: mohamadhasan.bahari@esat.kuleuven.be
(Mohamad Hasan Bahari), m.mclaren@let.ru.nl (Mitchell McLaren),
hugo.vanhamme@esat.kuleuven.be (Hugo Van hamme), d.vanleeuwen@let.ru.nl
(David A. van Leeuwen)

effect of some major factors influencing the proposed age estimation system, namely utterance length and spoken language are analyzed.

Keywords: speaker age estimation, i-vector, least squares support vector regression, utterance length, language mismatch.

1. Introduction

Automatic identification of age from speech signals has a wide range of commercial and forensic applications (Dobry et al., 2011; Tanner and Tanner, 2004; Li et al., 2013). For example, in targeted advertising through internet, where user-computer and user-company vocal interaction has increased significantly during the last decades, information about the user’s language/accent, age and gender can help to offer appropriate products and services (Schuller et al., 2013). Speaker age estimation is also required in many forensic scenarios such as kidnapping, threatening calls and false alarms to facilitate the identification of criminals, e.g. to narrow down the number of suspects (Tanner and Tanner, 2004). This technology can also guide ambient assisted living and smart home systems to automatically adapt to different user needs (Li et al., 2013). Speaker age estimation is also required in natural human-machine interaction. In video games, knowledge about the user’s age and gender can help to adapt the game accordingly (Schuller et al., 2013). Automatic identification of speaker age can be applied to improve the performance of other speech technology systems such as emotional state recognition, smoker speaker detection, identifying the level of intoxication and even automatic speech recognition (ASR).

Experimental studies reveal major effects of vocal aging on the speech sig-

21 nal such as lowered speaking rate and increased jitter and shimmer (Schotz,
 22 2006), and has shown to negatively influence speaker recognition perfor-
 23 mance (Kelly et al., 2013). However, the relation of these acoustic cues
 24 with speaker age is usually complex and affected by many other factors such
 25 as speech content, language, gender, weight, height, emotional condition,
 26 smoking and drinking habits (Schotz, 2006; Bahari and Van hamme, 2011;
 27 Bahari et al., 2012b). Furthermore, in many practical cases we have no con-
 28 trol over the available speech duration, content, language, etc.. These issues
 29 make automatic speaker age estimation very challenging for both humans
 30 and machines (Bocklet et al., 2008b; Li et al., 2013; Schotz, 2006).

31 Figure 1, which shows a simplified model for human speech production,
 32 helps to display the underling difficulties in speaker age estimation. In this
 33 problem, the recorded speech signal is the only available information and
 34 the task is to estimate one of the physical states of the articulatory system,
 35 namely the speaker’s age, without any information about the system inputs,
 36 channel characteristics and the other psychological and physical states of the
 37 articulatory system such as gender, emotional state and smoking habit.

38 Technical factors such as available speech duration, environment, record-
 39 ing device and channel conditions also influence the estimation accuracy. In
 40 other words, in a typical practical scenario, the quality of the available speech

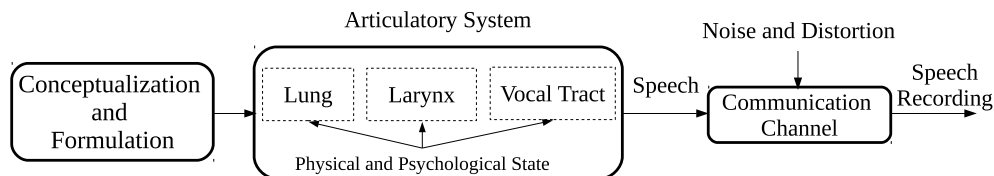


Figure 1: *simplified human speech production model and recording channel.*

41 signal and the recording conditions are not controlled and the duration of
42 the speech signal may vary from a few seconds to several hours.

43 *1.1. Related Work*

44 Studies on the influence of ageing on voice started in the late 1950s (Mysak,
45 1959). However, the first automatic speaker age recognition systems were de-
46 veloped around four decades later in the early 2000s (Linville, 2001; Muller
47 et al., 2003; Minematsu et al., 2002; Shafran et al., 2003). During this decade,
48 many different techniques, mostly inspired from the automatic speaker and
49 language recognition fields, have been suggested for categorizing speakers
50 based on their age groups. For example, using different types of acoustic
51 features and Support Vector Machines (SVM) (Mahmoodi et al., 2011; Chen
52 et al., 2011; van Heerden et al., 2010), Gaussian Mixture Model (GMM)
53 mean supervectors and SVM (Bocklet et al., 2008a), nuisance attribute pro-
54 jection (Dobry et al., 2009), anchor models (Dobry et al., 2009) and parallel
55 phoneme recognizers (Metze et al., 2007). The age sub-challenge of the Inter-
56 speech 2010 paralinguistic challenge provided a forum for presenting state-
57 of-the-art methods in speaker age group classification (Schuller et al., 2010).
58 Participants of the age sub-challenge tried to categorize speakers of telephony
59 data in the “aGender” corpus into four age groups — 7 to 14 (Child), 15 to
60 24 (Youth), 25 to 54 (Adult) and 55 to 80 (Senior) years old. In this sub-
61 challenge, GMM mean supervectors (Bocklet et al., 2010), GMM weight su-
62 pervectors (Porat et al., 2010), Maximum-Mutual-Information (MMI) train-
63 ing (Kockmann et al., 2010) and fuzzy SVM modeling (Nguyen et al., 2010)
64 have been suggested to enhance acoustic modeling quality. A brief overview
65 of different proposed methods in this sub-challenge is presented in (Li et al.,

2013), which also introduces an age group recognition approach using acoustic and prosodic level information fusion.

In speaker age group recognition, crisp borders are assumed between different age groups. For example, in the mentioned age sub-challenge, a speaker with age 54 belongs to the adult group and a 55 year old speaker belongs to the senior category. These two speakers who have only one year of age difference and share many similarities are considered to be from two different categories, while a 80 year old speaker with distant characteristics is in the same category as the 55 year old speaker. This setup causes many problems in training, testing, and performance measurement. To avoid these troubles, recently it has been suggested to use regression for age estimation (Bocklet et al., 2008b; Dobry et al., 2011; Bahari and Van hamme, 2011; Bahari et al., 2012b; Feld et al., 2009). A probabilistic interpretation of the posterior distribution of age estimation and its calibration is presented in (van Leeuwen and Bahari, 2012).

1.2. Motivations, Goals and Summary of Contributions

One effective approach to age estimation from speech involves modeling speech recordings with Gaussian Mixture Model (GMM) mean supervectors to use them as features in Support Vector Regression (SVR) (Dobry et al., 2011; Bocklet et al., 2008b). Similar Support Vector Machine (SVM) techniques have been successfully applied to different speech processing tasks such as speaker recognition (Campbell et al., 2006). While effective, GMM mean supervectors are of a high dimensionality resulting in high computational cost and difficulty in obtaining a robust model in the context of limited data. Consequently, dimension reduction through PCA-based methods has

91 been found to improve performance in age estimation from GMM mean su-
 92 pervectors (Dobry et al., 2011). In the field of speaker recognition, recent
 93 advances using so-called i-vectors (Dehak et al., 2011a) have increased the
 94 classification accuracy considerably. An i-vector is a compact representa-
 95 tion of an utterance in the form of a low-dimensional feature vector. The
 96 same idea was also applied in speaker or language and dialect recognition
 97 effectively (Dehak et al., 2011b; Bahari et al., 2013). In our last paper, we
 98 successfully replaced GMM mean supervectors by low-dimensional i-vectors
 99 to model utterances in an SVR based speaker age estimation system (Bahari
 100 et al., 2012a). The results of evaluation on the NIST 2010 and 2008 SRE
 101 databases illustrated that the i-vector based speaker age estimator increases
 102 the estimation accuracy.

103 In this paper, we extended our previous work (Bahari et al., 2012a) by:

- 104 1. Applying Within Class Covariance Normalization (WCCN) (Hatch et al.,
 105 2006) for session variability compensation. In our last paper (Bahari et al.,
 106 2012a), we have applied WCCN to normalize utterances of each age group.
 107 This method was not successful. In this paper we updated our strategy to
 108 use WCCN, where the classes are speakers rather than age groups.
- 109 2. Replacement of SVR by least squares SVR (LSSVR) to improve the com-
 110 putational cost.
- 111 3. Updating the evaluation setup to increase the size of training dataset,
 112 which helps the classifier to observe more variability in the data.
- 113 4. Using a standard z-test to analyze the statistical significance of the ob-
 114 tained improvements by the proposed method.
- 115 5. Investigate the effect of utterance length on the proposed automatic

116 speaker age estimation system.

117 6. Investigate the language mismatch on the proposed method.

118

119 The rest of this paper is organized as follows. In Section 2 the problem
120 of speaker age estimation and different conventional approaches addressing
121 this issue are described. In section 3, the proposed approach is elaborated.
122 Section 4 explains our experimental setup. The evaluation results and an
123 investigation of parameters affecting speaker age estimation are presented
124 and discussed in section 5. The paper ends with conclusions in section 6.

125 2. Age Estimation from Speech

126 In speaker age estimation, we are given a training dataset of speech
127 recordings $S^{\text{tr}} = \{(\mathcal{X}_1, y_1), \dots, (\mathcal{X}_s, y_s), \dots, (\mathcal{X}_S, y_S)\}$. In this set, \mathcal{X}_s and
128 y_s denote the s^{th} utterance of the training dataset and its corresponding
129 speaker age, respectively. The goal is to design an estimator function \mathcal{G} ,
130 such that for an utterance of an unseen speaker \mathcal{X}^{tst} , the actual speaker age
131 is predicted accurately.

132 2.1. Baseline Approaches

133 In this paper, we use three baseline approaches with which we compare
134 our proposed regression techniques:

135 **Prior:** The most basic choice for the estimator function is the average age
136 of the training data, $g(\mathcal{X}^{\text{tst}}) = \frac{1}{S} \sum_s y_s$. This estimator, labeled as *prior* in
137 the rest of this paper, intuitively provides a reference level of accuracy.

138 **GMM-R:** Different methods have been introduced to reach an effective
139 speaker age estimation (Dobry et al., 2011)–(Bahari and Van hamme, 2011).

For example, Bocklet et al. introduced GMM-R to estimate the age of children from GMM mean supervectors derived from their utterances (Bocklet et al., 2008b). Given an utterance, Maximum A Posteriori adaptation (MAP) is applied to adapt a Universal Background Model (UBM) to the speech characteristics of the speaker (Campbell et al., 2006). The component means of the obtained GMM are then extracted and concatenated to form a GMM mean supervector representing the utterance. Finally, an SVR is applied as a function approximator to estimate the speakers' age.

GMM-PCA-R and **GMM-WPPCA-R**: The approach of GMM-R was adopted and extended by Dobry et al. (Dobry et al., 2011) by applying dimension reduction techniques to the supervector. Methods such as Principal Component Analysis (PCA) and Weighted-Pairwise PCA (WPPCA) were applied and investigated. It was concluded that WPPCA, which is a supervised dimensionality reduction approach based on nuisance attribute projection (Dobry et al., 2011), yields more accurate results. These speaker age estimators, labeled GMM-PCA-R and GMM-WPPCA-R, are used as contrastive baseline systems in this paper.

3. Age Estimation using i-vectors

This section briefly describes the main components of the i-vector based age estimation approach, namely SVR and LSSVR, the i-vector framework and WCCN. Then, the proposed method is elaborated and finally the proposed scheme is presented.

163 3.1. Regression

164 In this section, SVR and LSSVR are briefly introduced.

165 3.1.1. Support Vector Regression

166 Support Vector Regression (SVR) is a function approximation approach
167 developed as a regression version of the widely known classification paradigm,
168 namely Support Vector Machines (SVM) (Lu et al., 2009; Smola and
169 Scholkopf, 2004). While SVMs perform the classification task by determin-
170 ing the maximum margin separation hyperplane between two classes, SVRs
171 carry out the regression task by finding the optimal regression hyperplane
172 in which most of training samples lie within an ϵ -margin around this hyper-
173 plane (Smola and Scholkopf, 2004). In a typical regression problem a training
174 dataset $S^{\text{tr}} = \{(a_1, b_1), \dots, (a_n, b_n), \dots, (a_N, b_N)\} \subset \mathbb{R}^d \times \mathbb{R}$ is given, where
175 a_n and b_n denote model input and corresponding output of the n^{th} data point
176 respectively. The objective of the regression analysis is to determine a func-
177 tion $f(a)$, so as to predict the desired outputs accurately. In the primal form
178 of SVR the following relation is considered for $f(a)$:

$$f(a) = \mathbf{w}^t \Phi(a) + z \quad (1)$$

179 where $\Phi(a)$ denotes a mapping function in the feature space, \mathbf{w} is a row vector
180 with the same dimension of $\Phi(a)$, $z \in \mathbb{R}$ is a constant and t represents
181 the transpose operator. Using Vapnik's ϵ -insensitive loss function the model
182 training—estimation of \mathbf{w} and z —is formulated as to minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{n=1}^N (\xi_n + \xi_n^*) \quad (2)$$

183 subject to

$$\begin{cases} b_n - \mathbf{w}^t \Phi(a_n) - z \leq \epsilon + \xi_n \\ \mathbf{w}^t \Phi(a_n) + z - b_n \leq \epsilon + \xi_n^* , \\ \xi_n, \xi_n^* \geq 0. \end{cases} \quad (3)$$

184 where ξ_n and ξ_n^* are slack variables vanishing during the optimization process,
 185 $\epsilon > 0$ controls the ϵ -insensitive zone used for fitting the training data and
 186 $\lambda > 0$ determines the trade-off between the atness of $f(a)$ and the cost of
 187 tolerating deviations larger than ϵ .

188 For high dimensional data, this constrained minimization problem can be
 189 solved more efficiently by introducing a dual set of variables and solving the
 190 following dual optimization problem (Smola and Scholkopf, 2004)

$$\begin{aligned} \max_{\alpha, \alpha^*} & -\frac{1}{2} \sum_{m,n=1}^N (\alpha_n - \alpha_n^*)(\alpha_m - \alpha_m^*) \langle \Phi(a_n), \Phi(a) \rangle \\ & - \epsilon \sum_{n=1}^N (\alpha_n - \alpha_n^*) + \sum_{n=1}^N (\alpha_n - \alpha_n^*) b_n, \end{aligned} \quad (4)$$

191 subject to the constraints

$$\begin{cases} \sum_{n=1}^N (\alpha_n - \alpha_n^*) = 0 \\ 0 \leq \alpha_n \leq \lambda, \quad n = 1, \dots, N, \\ 0 \leq \alpha_n^* \leq \lambda, \quad n = 1, \dots, N \end{cases} \quad (5)$$

192 where $\langle \cdot, \cdot \rangle$ describes the dot product and α and α^* are the dual set of vari-

ables. The resulting SVR model is

$$f(a) = \sum_{n=1}^N \beta_n \langle \Phi(a_n), \Phi(a) \rangle + z \quad (6)$$

$$= \sum_{n=1}^N \beta_n K(a_n, a) + z, \quad (7)$$

where $K(a_n, a)$ is the kernel function. Any function meeting the Mercer's condition can be used as the kernel function (Lu et al., 2009; Smola and Scholkopf, 2004). Parameters $\beta_n = \alpha_n - \alpha_n^*$ are calculated through solving the dual optimization problem and have the following relation to \mathbf{w}

$$\mathbf{w} = \sum_{n=1}^N \beta_n \Phi(a_n). \quad (8)$$

Since both the primal and dual optimization problem are convex, a unique optimal solution can be found efficiently using numerical methods such as quadratic programming (QP) (Smola and Scholkopf, 2004). Computing parameters β_n and z is explained in (Smola and Scholkopf, 2004) in detail.

In the baseline systems GMM-PCA-R and GMM-WPPCA-R (Dobry et al., 2011), SVR model training and testing is implemented using LIB-SVM (Chang and Lin, 2011) and the hyperparameters of the SVR such as the minimal error margin ϵ and error cost factor λ are tuned using the N -fold cross validation technique on the training dataset. In this research, we use the same toolbox and apply the same approach to tune the hyperparameters.

3.1.2. Least Squares Support Vector Regression

Least Squares Support Vector Machine (LSSVM), which is a variant of SVM, was introduced by Suykens and Vandewalle Suykens et al. (2002). It is

211 employed as a machine learning tool for classification, clustering and regres-
 212 sion tasks. Compared to SVM, LSSVM benefits from a faster training process
 213 because the quadratic programming problem of SVM is reduced to that of
 214 solving a system of linear equations. Furthermore, the LSSVM formulation
 215 involves fewer tuning parameters (Fodor, 2003). A continuous function can
 216 be fitted to the training data with a Least Squares Support Vector Regres-
 217 sor (LSSVR), a technique which shares many of the advantages of LSSVM
 218 classification. In the primal form of LSSVR, which is the same as SVR, the
 219 following relation is considered for $f(a)$

$$f(a) = \mathbf{w}^t \Phi(a) + z. \quad (9)$$

220 In LSSVR, a least squares loss function is applied instead of Vapnik's ϵ -
 221 insensitive loss function to simplify the formulations to minimize

$$\frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \gamma \sum_{n=1}^N e_n^2 \quad (10)$$

222 subject to

$$b_n = \mathbf{w}^t \Phi(a_n) + z + e_n, \quad (11)$$

223 where γ is a error cost factor playing the same role of λ in the SVR formu-
 224 lation and $e_n \in \mathbb{R}$ are error variables.

225 Similar to SVR, for high dimensional data this optimization problem can
 226 be solved more efficiently by introducing the Lagrangian variables ν and
 227 solving the following dual optimization problem (Suykens et al., 2002)

$$\mathbf{L}(\mathbf{w}, z, e, \nu) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{2} \gamma \sum_{n=1}^N e_n^2 \quad (12)$$

$$- \sum_{n=1}^N \nu_n \{ \mathbf{w}^t \Phi(a_n) + z + e_n - b_n \}. \quad (13)$$

One can solve this optimization problem directly by taking the partial derivative of \mathbf{L} with respect to \mathbf{w} , z , e and ν and setting the results to zero which leads to solving a linear system of equations. Inserting the obtained results to 9 leads to the regression function

$$f(a) = \sum_{n=1}^N \nu_n \langle \Phi(a_n), \Phi(a) \rangle + z \quad (14)$$

$$= \sum_{n=1}^N \nu_n K(a_n, a) + z, \quad (15)$$

where $K(a_n, a)$ is the kernel function and ν and z are the solution to optimization problem (12).

LSSVR has two advantages and one drawback compared to SVR. The first advantage of LSSVR is that its model training is faster as its dual form corresponds to solving a linear system which involves less computation time compared to a QP problem of SVR. The second advantage is that the LSSVR is faster to tune as its formulation involves fewer hyperparameters to tune (the minimal error margin ϵ is not used here). A drawback of this simplification is the loss of sparseness (ν is less sparse compared to β), which has been highlighted in literature (Suykens et al., 2000; Li et al., 2006).

In this research, the LSSVR models training and testing is implemented using LSSVMLab (Suykens et al., 2002) and the hyperparameters of the LSSVR are tuned on the training set using the N -fold cross validation technique.

3.2. The i -vector framework

The age estimation approaches described in section 2.1 are based on GMM mean supervectors and have been shown to yield reasonable performance. In

the related field of speaker recognition, GMM supervectors are commonplace. Recent progress in this field, however, has found an alternate method of modeling GMM supervectors that provides far superior speaker recognition performance (Dehak et al., 2011a). This technique, referred to as the i-vector framework, assumes the GMM mean supervector, \mathbf{M} , can be decomposed as

$$\mathbf{M} = \mathbf{u} + \mathbf{T}\mathbf{v} \quad (16)$$

where \mathbf{u} is the mean supervector of the UBM, \mathbf{T} spans a low-dimensional subspace (400 dimensions in this work) and subspace vector \mathbf{v} is treated as a latent variable with the standard normal prior and the i-vector is its maximum-a-posteriori (MAP) point estimate.

The subspace matrix \mathbf{T} is estimated via maximum likelihood in a large training dataset. An efficient procedure for training \mathbf{T} and MAP adaptation of i-vectors \mathbf{v} can be found in (Kenny et al., 2008). In this approach, i-vectors are the low-dimensional representation of an audio recording that can be used for classification and regression purposes.

3.3. *i-vector Session Compensation*

Session compensation is one of the most dominant topics in the speaker recognition field (McLaren and van Leeuwen, 2012; Dehak et al., 2011a). The main reason for using session compensation techniques is removing different session variabilities from the feature vectors (such as GMM supervectors or i-vectors) to allow the subsequent modeling approaches to better observe important between-class information. In this paper, we use Within-Class Covariance Normalization (WCCN) to normalize the within-class covariance of the i-vector space to the identity matrix (Hatch et al., 2006). In doing

so, directions of relatively high within-class variation will be attenuated and thus prevented from dominating the space (Hatch et al., 2006). The WCCN transformation matrix \mathbf{B}_W is found through Cholesky decomposition of

$$\left[\frac{1}{j} \sum_{j=1}^j \frac{1}{N_j} \sum_{i=1}^{N_j} (\mathbf{v}_j^i - \bar{\mathbf{v}}_j) (\mathbf{v}_j^i - \bar{\mathbf{v}}_j)' \right]^{-1} = \mathbf{B}_W \mathbf{B}_W', \quad (17)$$

where \mathbf{v}_j^i is the i^{th} i-vector in the j^{th} speaker, $\bar{\mathbf{v}}_j = \frac{1}{N_j} \sum_i^{N_j} \mathbf{v}_j^i$ is the mean of the observations for the j^{th} speaker, N_j denotes the number of utterances of the j^{th} speaker and j is the total number of speakers in the training dataset.

3.4. Train and Test

The principle of the proposed age estimation approach is illustrated in Figure 2. As it can be interpreted from this figure, in the training phase, each utterance in the training dataset is converted to an i-vector. Then, WCCN is used to remove the session variability as described in Section 3.3. Finally, the obtained vectors along with their corresponding chronological speaker age are used to train the regressor. In the testing phase, an i-vector is extracted from the utterance of an unseen speaker. Then, WCCN is used to remove the session variability. Finally, the trained regressor uses the obtained vector to estimate the chronological age of the test speaker.

The use of i-vectors for age estimation has several distinct advantages over GMM supervectors. Firstly, the relatively low dimensionality of i-vectors (400) significantly reduces the computational burden of model training and estimation compared to a GMM supervector dimensionality of greater than 12,000 used in this work. Secondly, subspace adaptation of i-vector \mathbf{v} results

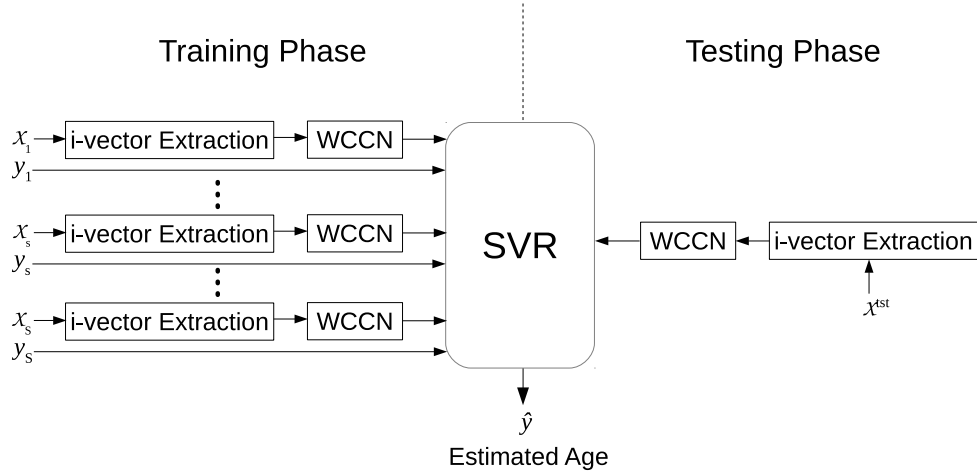


Figure 2: The block diagram of the proposed speaker age estimation approach in training and testing phases.

in a more reliable estimation of the true model means \mathbf{M} in the context of limited training data (Dehak et al., 2011b).

4. Experimental Setup

4.1. Database

The National Institute of Standards and Technology (NIST) has held annual or biannual Speaker Recognition Evaluations (SRE) for the past two decades. With each SRE, a large database of telephone conversations (and more recently microphone speech) are released along with an evaluation protocol. These conversations typically last five minutes and originate from a large number of participants for whom meta data is recorded—including participant age and language. The NIST databases were chosen for this work due to the large number of speakers meeting the i-vector framework

305 requirement for a considerable amount of development data to estimate sub-
 306 space matrix \mathbf{T} accurately. In our experiments, first a development dataset
 307 is formed, which includes over 30,000 speech recordings sourced from NIST
 308 2004–2006 SRE databases, to estimate the parameters of UBM and the sub-
 309 space matrix (\mathbf{T}). The procedure of obtaining the applied UBM and subspace
 310 matrix is presented in (McLaren and van Leeuwen, 2012).

311 To form the train and test datasets for speaker age estimation, telephone
 312 recordings from the common protocols of the NIST 2010 and 2008 SRE cor-
 313 pora are used. The core protocol, short2-short3, from the 2008 database
 314 contains 3772 telephone recordings from 1154 speakers for whom the age
 315 is between 20 and 70. The language label of 3726 utterances is given in
 316 this database. Among these, 2656 utterances are English and the remaining
 317 1070 utterances are from 26 different non-English languages including Rus-
 318 sian, Italian and Japanese. Similarly, the extended core-core protocol of the
 319 2010 database contains 5479 telephone speech segments from 422 speakers
 320 for whom the age is between 20 and 70. All utterances of this database are
 321 English. There is no overlap between speech recordings extracted from the
 322 NIST 2010 and NIST 2008 SRE databases.

323 Figure 3 illustrates the age histograms of male and female speakers in
 324 the NIST 2010 and 2008 SRE databases. Since the perceptions of speaker
 325 gender and age have a significant mutual impact, all the experiments are
 326 performed for male and female speakers separately in this paper.

327 4.2. Performance Metric

328 The effectiveness of the applied methods is evaluated using the Mean
 329 Absolute Error (E_{ma}) of the estimated speakers’ age and Pearson’s correlation

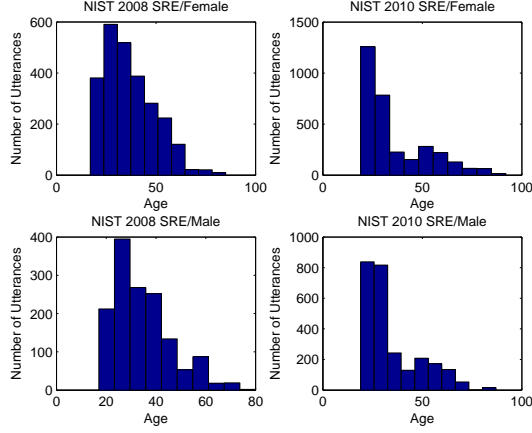


Figure 3: *Age histogram of telephone speech utterances for NIST 2010 and 2008 SRE Databases.*

330 coefficient (ρ) between chronological speakers' age and estimated speakers'
 331 age. The measure E_{ma} is calculated using:

$$E_{\text{ma}} = \frac{1}{Q} \sum_{q=1}^Q |\hat{y}_q - y_q|, \quad (18)$$

332 where \hat{y}_q and y_q are the estimated and the chronological age of the q^{th} utter-
 333 ance of the testing dataset respectively. Q is the total number of utterances
 334 in the testing dataset. Further,

$$\rho = \frac{1}{Q-1} \sum_{q=1}^Q \left(\frac{\hat{y}_q - \mu_{\hat{y}}}{\sigma_{\hat{y}}} \right) \left(\frac{y_q - \mu_y}{\sigma_y} \right), \quad (19)$$

335 where $\mu_{\hat{y}}$ and $\sigma_{\hat{y}}$ are the mean and the standard deviation of the speakers'
 336 estimated age respectively. Similarly μ_y and σ_y denote the mean and the
 337 standard deviation of the speakers' chronological age respectively.

338 We also apply the standard z-test to analyze the statistical significance
 339 level of differences between the mean absolute errors of applied systems.

340 5. Results and Discussion

341 This section presents the evaluation results of the baseline systems and
342 compares them to the introduced i-vector based age estimation system.

343 The applied GMM in all experiments consist of 512 mixture components.
344 To study the effect of the acoustic features, two types of feature vectors
345 have been tested for the baseline systems. The first type, labeled MFCC_{26D},
346 consists of 13 Mel-Frequency Cepstrum Coefficients (MFCCs) including ap-
347 pended energy with their first order derivatives, forming a 26 dimensional
348 acoustic feature vector. The second type, MFCC_{60D}, consists of 20 MFCCs
349 including appended energy with their first and second order derivatives, form-
350 ing a 60 dimensional acoustic feature vector. In both cases, a hamming win-
351 dow is used and the sampling rate, frame rate, frame size and number of Mel
352 frequency channels are 8000 Hz, 100 Hz, 0.02 s and 30 respectively. To have
353 more reliable features, Wiener filtering, speech activity detection (McLaren
354 and van Leeuwen, 2011) and feature warping (Pelecanos and Sridharan, 2001)
355 have been applied as front-end processing. The former type, MFCC_{26D},
356 matches the configuration of features applied in (Dobry et al., 2011) and
357 the latter type, MFCC_{60D}, is very common in state-of-the-art i-vector based
358 speaker recognition systems.

359 5.1. SVR and LSSVR

360 In this section, an experiment is performed to investigate the perfor-
361 mances of SVR and LSSVR for regression in this problem and to choose the
362 regression method with more accurate estimation results for the rest of the
363 experiments in this paper.

364 In this experiment, the NIST 2008 and 2010 SRE databases are used
 365 for training and testing respectively and the acoustic features are MFCC_{26D}.
 366 Each utterance in the training and testing datasets is modeled using its
 367 corresponding GMM mean supervector. Then, an SVR or an LSSVR are
 368 applied as a function approximator to estimate the speakers' age.

369 Like the baseline systems GMM-PCA-R and GMM-WPPCA-R, SVR
 370 model training and testing is performed using LIBSVM (Chang and Lin,
 371 2011) and the SVR Hyperparameters ϵ and λ are tuned using the 5-fold
 372 cross-validation. Since it is shown in (Dobry et al., 2011) that the radial
 373 basis function (RBF) kernel leads to more accurate estimation compared to
 374 the linear kernel, we apply the RBF kernel in our experiments. Two methods
 375 are applied to determine the width of the Gaussian functions. In the first
 376 scheme, which is adopted from (Dobry et al., 2011), the width of the Gaus-
 377 sian functions was set to $\sqrt{\det(\Sigma)/2}$, where Σ is the training feature vectors
 378 covariance matrix and $\det(\cdot)$ denotes the determinant operator. It was men-
 379 tioned in (Dobry et al., 2011) that $\sqrt{\det(\Sigma)/2}$ was found to be optimal on a
 380 number of empirical experiments. The results of this method, labeled as SVR
 381 1, are listed in the first row of Table 1. In the second approach, labeled as
 382 SVR 2, the 5-fold cross-validation is used to tune the width of the Gaussian
 383 functions.

384 The LSSVR approach applied in this experiment also uses the RBF kernel
 385 and 5-fold cross-validation in order to tune its error cost factor and Gaussian
 386 width.

387 Table 1 shows the obtained results using SVR 1, SVR 2 and LSSVR in
 388 this experiment. This table shows that LSSVR estimates the speakers' age

Table 1: *The E_{ma} (in years) and ρ of male and female speakers’ age estimation using SVR and LSSVR.*

Regression	Female		Male	
Method	E_{ma}	ρ	E_{ma}	ρ
SVR 1	7.59	0.80	7.97	0.69
SVR 2	7.48	0.80	7.92	0.70
LSSVR	7.44	0.80	7.87	0.70

more accurately compared to SVR 1 and SVR 2 in this experiment. LSSVR is selected for the rest of experiments in this paper rather than conventional SVR due to the obtained marginal improvement and faster and easier model training and tuning.

5.2. Baseline Systems Results

In this section, the performances of baseline systems, namely prior, GMM-R, GMM-PCA-R and GMM-WPPCA-R, are investigated.

To evaluate the baseline systems on all available utterances, 15-fold cross-validation is used. Therefore, first all speakers in the NIST 2008 and 2010 SRE databases are divided into 15 disjoint folds. Then, 15 independent experiments are run so that in each experiment, a new fold is used as the testing dataset and the remaining 14 folds are used as training dataset. Due to high variability in our data such as language, smoking habit and content, in our experiments, we have applied 15-fold rather than 5-fold or 10-fold cross-validation to have larger training datasets, which include more variability of the data.

405 The average E_{ma} and ρ of male and female speakers' age estimation us-
 406 ing the baseline systems in all 15 experiments with both types of acoustic
 407 features are listed in tables 2 and 3 respectively. In this experiment, PCA
 408 and WPPCA have been tested over different target dimensions between 100
 409 and 1000. Tables 2 and 3 only include the best results, which were obtained
 410 for target dimensions 300 and 400 for GMM-PCA-R and GMM-WPPCA-R
 411 respectively.

412 Results in tables 2 and 3 indicate that the GMM-R system is remarkably
 413 more accurate than the prior system. This shows that the GMM supervectors
 414 contain speaker information including age. The Tables 2 and 3 also show that
 415 the PCA and WPPCA based systems outperform the GMM-R system, thus
 416 demonstrating the benefit of dimension reduction of the GMM supervectors
 417 prior to regression. Unlike (Dobry et al., 2011) our experiments do not show
 418 remarkable advantage for using WPPCA over PCA. It is also interpreted from
 419 tables 2 and 3 that increasing the acoustic dimension from 26 to 60 slightly
 420 improves the estimation accuracy for GMM-PCA-R and GMM-WPPCA-R.
 421 Therefore, in the rest of our experiments we focused on the second type of
 422 acoustic features, MFCC_{60D}.

423 5.3. *i-vectors for Age Estimation*

424 The results of the proposed method for speakers' age estimation are pre-
 425 sented in this section.

426 Figures 4 and 5 present the E_{ma} of the estimated age and the ρ between
 427 the chronological speakers' age and the estimated speakers' age using the
 428 proposed method and the baseline systems for different target dimensions
 429 respectively. These figures show that the proposed method, labeled i-vector-

Table 2: The average E_{ma} (in years) of male and female speakers' age estimation for the baseline systems using MFCC_{26D} and MFCC_{60D} feature vectors.

System Configuration	Female		Male	
	MFCC _{26D}	MFCC _{60D}	MFCC _{26D}	MFCC _{60D}
Prior	10.57	10.57	10.08	10.08
GMM-R	6.19	6.60	6.93	7.53
GMM-PCA-R	6.26	6.21	6.79	6.71
GMM-WPPCA-R	6.25	6.17	6.74	6.74

WCCN-R, is more accurate than the other state-of-the-art approaches. Note that this improvement was obtained without any optimization over the target dimension in the i-vector framework. Therefore, in figures 4 and 5, the result of proposed method is only shown for dimension 400. In the standard i-vector framework, the optimization over the target dimension is usually very time-

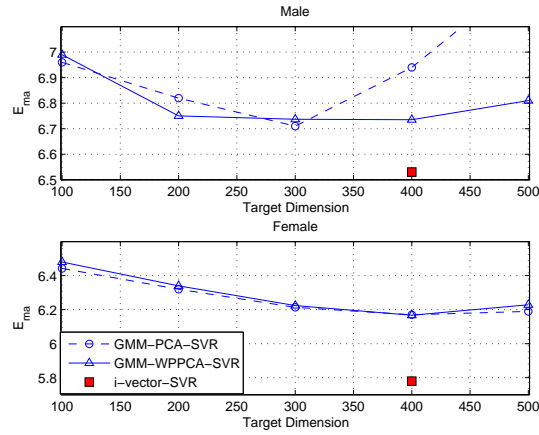


Figure 4: The E_{ma} of female and male speakers' age estimation using the proposed method and baseline systems versus target dimension.

Table 3: The average ρ of male and female speakers' age estimation for the baseline systems using MFCC_{26D} and MFCC_{60D} feature vectors.

System Configuration	Female		Male	
	MFCC _{26D}	MFCC _{60D}	MFCC _{26D}	MFCC _{60D}
Prior	0	0	0	0
GMM-R	0.78	0.73	0.69	0.59
GMM-PCA-R	0.77	0.78	0.71	0.72
GMM-WPPCA-R	0.77	0.78	0.71	0.71

435 consuming and computationally expensive. Furthermore, different studies
 436 such as (Dehak et al., 2011b) show that i-vector characteristics are mostly
 437 robust against different dimensions between 200 to 500.

438 The ρ and E_{ma} of age estimation using the proposed approach are 0.772
 439 and 6.08 respectively. Therefore, the proposed method improves ρ by 12.9%,

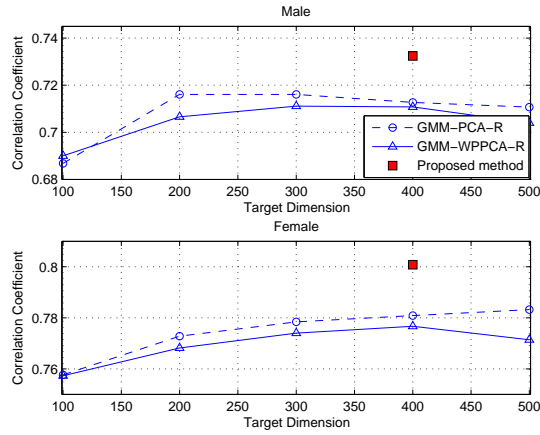


Figure 5: Pearson correlation coefficient between estimated and true age of female and male speakers using the proposed method and baseline systems versus target dimension.

2.0% and 2.6% relative to GMM-R, GMM-PCA-R and GMM-WPPCA-R respectively. The E_{ma} is also improved by 41%, 13%, 5% and 4.8% relative to Prior, GMM-R, GMM-PCA-R and GMM-WPPCA-R respectively. A standard z-test for comparing two means show that the E_{ma} of the i-vector based system method is significantly lower than that of the best baseline system, namely GMM-PCA-R, at the 99% confidence level. Details of this test are presented in Appendix A.

The coefficient of determination (ρ^2) (Draper and Smith, 1981) obtained by the proposed method for male and female speakers are 0.54 and 0.64 respectively, which show that roughly 54% and 64% of the variance in male and female speakers' age were successfully accounted for.

We also investigated using i-vectors without session variability compensation, like our earlier work (Bahari et al., 2012a). In this case, the ρ and E_{ma} are 0.76 and 6.22 respectively. This experiment shows that session variability compensation using WCCN relatively improves the ρ and E_{ma} by 1.5% and 2.2% respectively.

5.4. The Effect of Utterance Length

In a typical practical case, the duration of the available speech sample may vary from a few seconds to several hours. Although there is literature on the effect of available utterance duration on speaker recognition systems (Mandasari et al., 2011), there is no published research on this topic for automatic speaker age estimation systems. In this section, we analyze the performance of the proposed i-vector based speaker age estimation system with respect to speech duration in the terms of E_{ma} and ρ .

In this experiment, first all speakers in the NIST 2008 and 2010 SRE

Table 4: The E_{ma} and ρ of speakers’ age estimation using the proposed method in different test utterance length conditions.

Utterance Length	Female		Male	
	E_{ma}	ρ	E_{ma}	ρ
5s	9.51	0.53	8.99	0.47
10s	8.5	0.64	8.27	0.57
20s	7.38	0.72	7.66	0.64
45s	6.47	0.77	6.99	0.70

465 databases are divided into 15 disjoint folds. Then, 15 independent exper-
 466 iments are run so that in each experiment, a new fold is used as testing
 467 dataset and the rest 14 folds are used as training dataset. Each utterance in
 468 the testing dataset typically contains around 80 seconds of active speech. In
 469 order to study the effect of test sample duration, we synthesized test datasets
 470 of 5, 10, 20 and 40 seconds by truncating the feature streams after speech
 471 activity detection. For consistency in our results, the test samples that con-
 472 tained less than 40 seconds of nominal speech using our speech detection
 473 algorithm were discarded from all results reported in this experiment. The
 474 procedure and details of obtaining corresponding i-vectors for truncated test
 475 samples is explained in (Mandasari et al., 2013). The corresponding E_{ma}
 476 and ρ values are listed in Table 4. The performance of the proposed method
 477 decreases as the test utterance duration is reduced. This is more evident
 478 when the utterance duration is less than 10 seconds. However, the results
 479 of the proposed method remain significantly more accurate than the prior,
 480 even for the utterances with a length of 5 seconds.

481 5.5. *The Effect of Language*

482 Braun and Cerrato performed a number of experiments to evaluate the
 483 ability of human listeners in estimating speakers’ age across different lan-
 484 guages (Braun and Cerrato, 1999). They concluded that the age can be
 485 estimated almost as accurately when the listeners are familiar with the lan-
 486 guage of the speaker as when they are not. However, Schotz considered the
 487 language as an important source influencing the acoustic analysis of speaker
 488 age (Schotz, 2006). Feld et al. studied the effect of language mismatch be-
 489 tween train database and test samples on automatic speaker age estimation
 490 systems. In this section, we analyze the effect of language mismatch on the
 491 proposed i-vector based age estimation system.

492 In this experiment, the train database is NIST 2010 SRE, which includes
 493 5634 English utterances from 445 speakers. There are two test databases in
 494 this experiment, the English and non-English parts of the NIST 2008 SRE
 495 database. Figure 6 illustrates the age histograms of the English and non-
 496 English speakers of the NIST 2008 SRE database. To eliminate the effect
 497 of utterance length, we synthesized test samples of 40 seconds by truncating
 498 the feature streams after speech activity detection. The E_{ma} and ρ of this
 499 experiment for both English and non-English test sets are listed in table 5.

Table 5: *The E_{ma} and ρ for both English and non-English test sets.*

System Configuration	Female		Male	
	English	Non-English	English	Non-English
E_{ma}	6.92	8	7.72	8.32
ρ	0.66	0.42	0.50	0.32

Results in table 5 indicate that language mismatch between train database and test samples causes a large performance degradation in both E_{ma} and ρ . It is obvious that the E_{ma} for the English test set is significantly less than that of the non-English test set for both male and female utterances. In these experiments, since only telephone speech signals are used, we do not concentrate on channel mismatch. The effect of gender is also discarded because all the experiments are performed for male and female speakers separately.

6. Conclusions

In this paper, utterance modeling with i-vectors, which was successfully applied to speaker recognition, has been used in conjunction with a WCCN and a LSSVR to address speaker age estimation. For the evaluation, telephone utterances of NIST 2010 and 2008 SRE databases have been used.

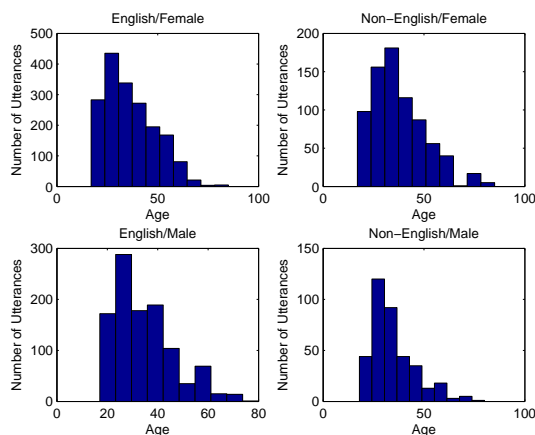


Figure 6: *Age histogram of English and non-English speakers in the NIST 2008 SRE database.*

513 Assessment results demonstrate that ρ and E_{ma} for the proposed approach
 514 are 0.772 and 6.08 respectively. Therefore, the obtained relative improve-
 515 ments of ρ and E_{ma} compared to the best baseline system are around 2%
 516 and 5% respectively. The experiments on analyzing the effect of utterance
 517 duration reveals that the performance of the proposed method degrades as
 518 the utterance length decreases especially for samples shorter than 20 sec-
 519 onds. However, it is still more accurate than the prior baseline system even
 520 for utterances of 5 seconds in length. Analyzing the effect of language shows
 521 that the language mismatch between train and test databases significantly
 522 decreases the performance of the age estimation system.

523 **7. Acknowledgements**

524 This work is supported by the European Commission as a Marie-Curie
 525 ITN-project (FP7-PEOPLE-ITN-2008), namely Bayesian Biometrics for
 526 Forensics (BBfor2), under Grant Agreement number 238803.

527 **Appendix A.**

528 In this appendix, a statistical analysis is presented to compare the mean
 529 absolute errors of age estimation obtained by the i-vector-SVR and GMM-
 530 PCA-R.

531 Since the values of populations variances are unknown, tests for the com-
 532 parison of two means should be conducted with the a t -test normally. How-
 533 ever, both sample sizes are greater than 30 in this case and we can work with
 534 the standard normal distribution (z -test) instead of Student distribution (t -
 535 test). In the standard z -test for comparison of two means, the z value is

536 calculated as follows:

$$z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (\text{A.1})$$

537 where \bar{x}_i , s_i , and n_i denote the mean, the variance and total number of
538 samples in the first set respectively. Similarly, \bar{x}_2 , s_2 , and n_2 are the mean,
539 the variance and sample size in the second set respectively.

540 In the comparison of the mean absolute errors of age estimation obtained
541 by the i-vector-SVR (\bar{x}_1) and GMM-PCA-R (\bar{x}_2), the null hypothesis is $\bar{x}_2 \leq$
542 \bar{x}_1 and the alternative hypothesis is $\bar{x}_2 > \bar{x}_1$. With significance levels $\alpha = 0.01$
543 and $\alpha = 0.05$, the critical values of z are 2.33 and 1.645 respectively for a
544 one tail test.

545 The mean and the standard deviation of the age estimation absolute error
546 using i-vector-SVR and GMM-PCA-R over male and female utterances are
547 listed in Table A.6.

548 As it is shown in Table A.6, the obtained z for male and female utterances
549 is greater than critical value of z for significance levels $\alpha = 0.05$ and $\alpha = 0.01$
550 respectively. Therefore, the null hypothesis is rejected and it is concluded
551 that the alternative hypothesis is true.

552 In the test of significance, we are trying to compare GMM-WPPCA-R
553 and the proposed method. Consequently, all results of the proposed method
554 (regardless of gender) can be considered in one class and all the results of
555 GMM-WPPCA-R are assumed to be in the other class. The last row of
556 Table A.6 shows the mean and the standard deviation of the age estimation
557 absolute error using the proposed method and GMM-WPPCA-R over all
558 utterances regardless of gender (labeled both). The obtained z value of this
559 experiment is 4.15 which is greater than the critical value of z for significance

Table A.6: *The mean and the standard deviation of age estimation absolute error using i-vector-SVR and GMM-PCA-R over male and female utterances.*

Gender	Parameter	Proposed method	GMM-WPPCA-R	z
Male	\bar{x}_i	6.53	6.74	1.7
	s_i	5.36	5.54	
	n_i	3883	3883	
Female	\bar{x}_i	5.78	6.17	4.13
	s_i	4.78	4.92	
	n_i	5292	5292	
Both	\bar{x}_i	6.10	6.41	4.15
	s_i	5.05	5.20	
	n_i	9175	9175	

level $\alpha = 0.01$.

References

- Bahari, M.H., McLaren, M., Van Leeuwen, D., Van hamme, H., 2012a. Age estimation from telephone speech using i-vectors, in: Proc. Interspeech, pp. 506–509.
- Bahari, M.H., Saeidi, R., Van hamme, H., van Leeuwen, D., 2013. Accent recognition using i-vector, gaussian mean supervector and gaussian posterior probability supervector for spontaneous telephone speech, in: International conference on acoustics, speech, and signal processing, pp. 7344–7348.

- 570 Bahari, M.H., Van hamme, H., 2011. Speaker age estimation and gender
571 detection based on supervised non-negative matrix factorization, in: Proc.
572 IEEE Workshop on Biometric Measurements and Systems for Security and
573 Medical Applications (BIOMS), pp. 1–6.
- 574 Bahari, M.H., et al., 2012b. Speaker age estimation using hidden markov
575 model weight supervectors, in: 11th IEEE Int. Conf. Information Science,
576 Signal Processing and their Applications (ISSPA), pp. 517–521.
- 577 Bocklet, T., Maier, A., Bauer, J., Burkhardt, F., Noth, E., 2008a. Age and
578 gender recognition for telephone applications based on GMM supervectors
579 and support vector machines, in: Proc. IEEE Int. Conf. Acoustics, Speech
580 and Signal Processing (ICASSP), pp. 1605–1608.
- 581 Bocklet, T., Maier, A., Noth, E., 2008b. Age determination of children
582 in preschool and primary school age with GMM-based supervectors and
583 support vector machines-regression, in: Proc. Text, Speech and Dialogue,
584 pp. 253–260.
- 585 Bocklet, T., Stemmer, G., Zeissler, V., Noth, E., 2010. Age and gender
586 recognition based on multiple systems early vs. late fusion, in: Proc. 11th
587 Annual Conference of the International Speech Communication Association,
588 pp. 2830–2833.
- 589 Braun, A., Cerrato, L., 1999. Estimating speaker age across languages, in:
590 Proc. ICPHS, pp. 1369–1372.
- 591 Campbell, W., Sturim, D., Reynolds, D., 2006. Support vector machines

592 using GMM supervectors for speaker verification. *IEEE Signal Processing*
593 *Letters* 13, 308–311.

594 Chang, C.C., Lin, C.J., 2011. Libsvm: a library for support vector machines.
595 *ACM Transactions on Intelligent Systems and Technology (TIST)* 2, 27.

596 Chen, C.C., Lu, P.T., Hsia, M.L., Ke, J.Y., Chen, O.C., 2011. Gender-to-age
597 hierarchical recognition for speech, in: *Circuits and Systems (MWSCAS),*
598 *2011 IEEE 54th International Midwest Symposium on*, pp. 1–4.

599 Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., Ouellet, P., 2011a. Front-
600 end factor analysis for speaker verification. *IEEE Trans. Audio, Speech,*
601 *and Lang. Process.* 19, 788–798.

602 Dehak, N., Torres-Carrasquillo, P.A., Reynolds, D., Dehak, R., 2011b. Lan-
603 guage recognition via ivectors and dimensionality reduction, in: *Proc. In-*
604 *terspeech*, pp. 857–860.

605 Dobry, G., Hecht, R., Avigal, M., Zigel, Y., 2009. Dimension reduction
606 approaches for svm based speaker age estimation, in: *Proc. Interspeech,*
607 pp. 2031–2034.

608 Dobry, G., Hecht, R.M., Avigal, M., Zigel, Y., 2011. Supervector dimension
609 reduction for efficient speaker age estimation based on the acoustic speech
610 signal. *IEEE Trans. Audio, Speech, and Lang. Process.* 19, 1975–1985.

611 Draper, N.R., Smith, H., 1981. *Applied regression analysis* 2nd ed. .

612 Feld, M., Barnard, E., van Heerden, C., Muller, C., 2009. Multilingual
613 speaker age recognition: Regression analyses on the lwazi corpus, in: *Au-*

614 automatic Speech Recognition & Understanding, 2009. ASRU 2009. IEEE
615 Workshop on, pp. 534–539.

616 Fodor, I., 2003. Statistical techniques to find similar objects in images, in:
617 Proc. the American Statistical Association, Statistical Computing Section.

618 Hatch, A., Kajarekar, S., Stolcke, A., 2006. Within-class covariance normal-
619 ization for svm-based speaker recognition, in: Proc. Interspeech.

620 van Heerden, C., Barnard, E., Davel, M., van der Walt, C., van Dyk, E., Feld,
621 M., Muller, C., 2010. Combining regression and classification methods for
622 improving automatic speaker age recognition, in: Acoustics Speech and
623 Signal Processing (ICASSP), 2010 IEEE International Conference on, pp.
624 5174–5177.

625 Kelly, F., Drygajlo, A., Harte, N., 2013. Speaker verification in score-ageing-
626 quality classification space. *Computer Speech & Language* .

627 Kenny, P., Ouellet, P., Dehak, N., Gupta, V., Dumouchel, P., 2008. A study
628 of interspeaker variability in speaker verification. *IEEE Trans. Audio,
629 Speech, and Lang. Process.* 16, 980–988.

630 Kockmann, M., Burget, L., Cernocky, J., 2010. Brno university of technol-
631 ogy system for interspeech 2010, in: Proc. 11th Annual Conference of the
632 International Speech Communication Association, pp. 2822–2825.

633 van Leeuwen, D.A., Bahari, M.H., 2012. Calibration of probabilistic age
634 recognition, in: Proc. Interspeech, pp. 502–505.

635 Li, M., Han, K.J., Narayanan, S., 2013. Automatic speaker age and gender
636 recognition using acoustic and prosodic level information fusion. *Computer*
637 *Speech and Language* 27, 151 – 167.

638 Li, Y., Lin, C., Zhang, W., 2006. Improved sparse least-squares support
639 vector machine classifiers. *Neurocomputing* 69, 1655–1658.

640 Linville, S.E., 2001. Vocal aging. Singular Thomson Learning.

641 Lu, C., Lee, T., Chiu, C., 2009. Financial time series forecasting using
642 independent component analysis and support vector regression. *Decision*
643 *Support Systems* 47, 115–125.

644 Mahmoodi, D., Soleimani, A., Marvi, H., Razzazi, F., Taghizadeh, M., Mah-
645 moodi, M., 2011. Age estimation based on speech features and support
646 vector machine, in: 3rd Computer Science and Electronic Engineering Con-
647 ference, pp. 60–64.

648 Mandasari, M., McLaren, M., van Leeuwen, D., 2011. Evaluation of i-vector
649 speaker recognition systems for forensic application, in: *Proc. Interspeech*,
650 pp. 21–24.

651 Mandasari, M.I., Saeidi, R., McLaren, M., van Leeuwen, D.A., 2013. Quality
652 measure functions for calibration of speaker recognition system in vari-
653 ous duration conditions. *IEEE Transactions on Acoustics, Speech, and*
654 *Language Processing* 21, 2425 – 2438.

655 McLaren, M., van Leeuwen, D., 2011. A simple and effective speech activity
656 detection algorithm for telephone and microphone speech. *Proc. NIST*
657 *SRE Workshop* , 1–6.

- 658 McLaren, M., van Leeuwen, D., 2012. Source-normalized lda for robust
 659 speaker recognition using i-vectors from multiple speech sources. IEEE
 660 Trans. Audio, Speech, and Lang. Process. 20, 755–766.
- 661 Metze, F., Ajmera, J., Englert, R., Bub, U., Burkhardt, F., Stegmann, J.,
 662 Muller, C., Huber, R., Andrassy, B., Bauer, J., et al., 2007. Comparison
 663 of four approaches to age and gender recognition for telephone applica-
 664 tions, in: Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing
 665 (ICASSP), pp. IV–1089.
- 666 Minematsu, N., Sekiguchi, M., Hirose, K., 2002. Automatic estimation of
 667 one’s age with his/her speech based upon acoustic modeling techniques of
 668 speakers, in: IEEE Int. Conf. Acoustics, Speech, and Signal Processing
 669 (ICASSP), pp. I–137.
- 670 Muller, C., Wittig, F., Baus, J., 2003. Exploiting speech for recognizing
 671 elderly users to respond to their special needs, in: Proc. 8th European
 672 Conf. Speech Communication and Technology (Eurospeech), pp. 1305–
 673 1308.
- 674 Mysak, E.D., 1959. Pitch and duration characteristics of older males. Journal
 675 of Speech, Language and Hearing Research 2, 46.
- 676 Nguyen, P., Le, T., Tran, D., Huang, X., Sharma, D., 2010. Fuzzy support
 677 vector machines for age and gender classification, in: Proc. 11th Annual
 678 Conference of the International Speech Communication Association, pp.
 679 2806–2809.

680 Pelecanos, J., Sridharan, S., 2001. Feature warping for robust speaker veri-
681 fication , 213–218.

682 Porat, R., Lange, D., Zigel, Y., 2010. Age recognition based on speech
683 signals using weights supervector, in: Proc. 11th Annual Conference of the
684 International Speech Communication Association, pp. 2814–2817.

685 Schotz, S., 2006. Perception, analysis and synthesis of speaker age. volume 47.
686 Citeseer.

687 Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Muller, C.,
688 Narayanan, S., 2010. The interspeech 2010 paralinguistic challenge, in:
689 Proc. Interspeech, pp. 2794–2797.

690 Schuller, B., Steidl, S., Batliner, A., Burkhardt, F., Devillers, L., Müller, C.,
691 Narayanan, S., 2013. Paralinguistics in speech and language state-of-the-art
692 and the challenge. *Computer Speech & Language* 27, 4–39.

693 Shafran, I., Riley, M., Mohri, M., 2003. Voice signatures, in: Proc. IEEE
694 Workshop on Automatic Speech Recognition and Understanding (ASRU),
695 pp. 31–36.

696 Smola, A., Scholkopf, B., 2004. A tutorial on support vector regression.
697 *Statistics and computing* 14, 199–222.

698 Suykens, J.A., Lukas, L., Vandewalle, J., 2000. Sparse approximation using
699 least squares support vector machines, in: *Circuits and Systems, 2000. Pro-*
700 *ceedings. ISCAS 2000 Geneva. The 2000 IEEE International Symposium*
701 *on*, pp. 757–760.

702 Suykens, J.A.K., et al., 2002. Least squares support vector machines. World
703 Scientific.

704 Tanner, D.C., Tanner, M.E., 2004. Forensic aspects of speech patterns: voice
705 prints, speaker profiling, lie and intoxication detection. Lawyers and Judges
706 Publishing.